# DEFECT INSPECTION FROM SCRATCH TO PRODUCTION

Dr. Andrew Liu

Sr. Solution Architect

**NVIDIA**

# AGENDA

Image Segmentation
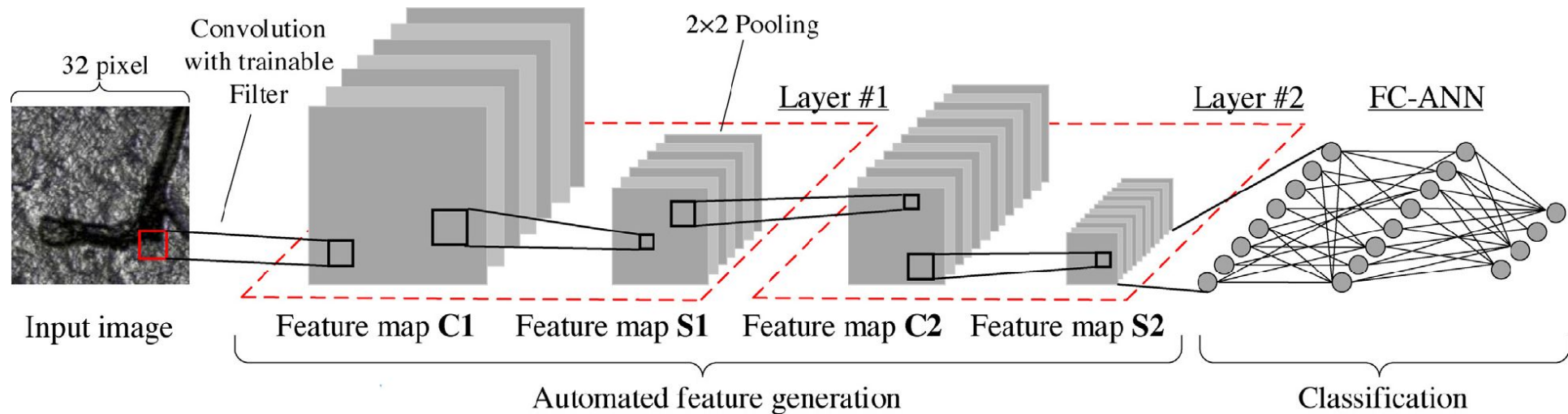
    - Fully Convolution Neural Network

Defect Inspection

    - Problem Define

    - Data Preparation

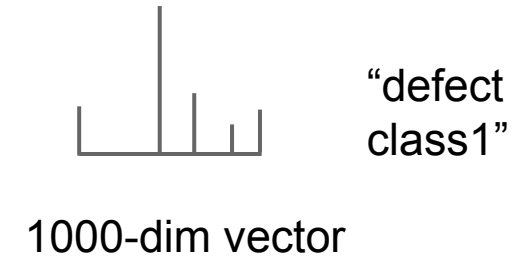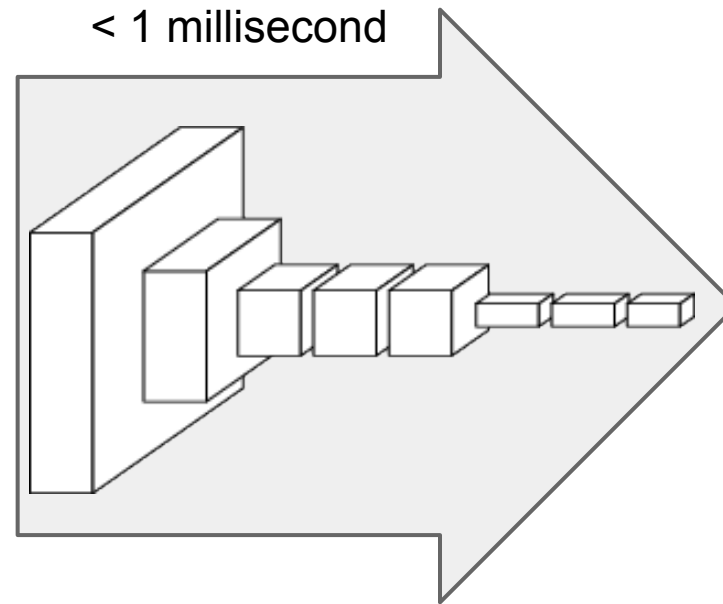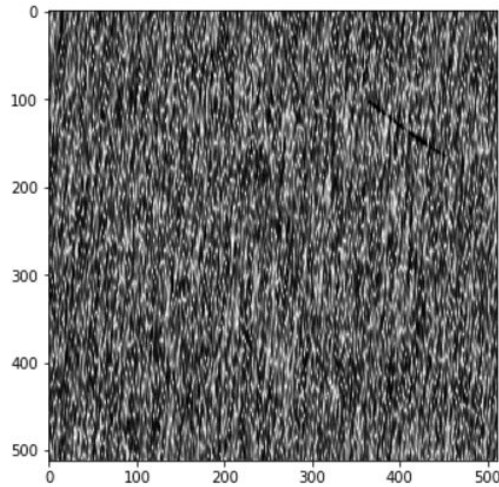    - Deal with Imbalance data
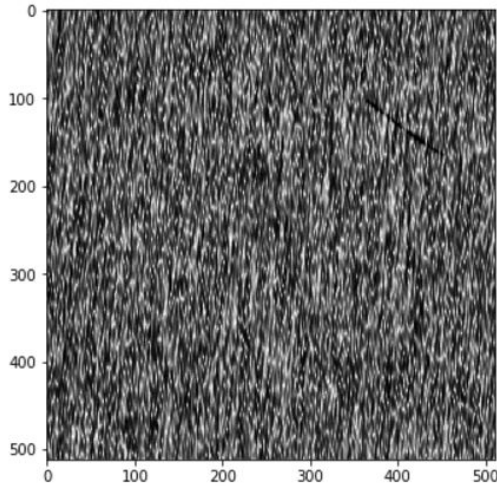
Speed up with TensorRT

# CNN STRUCTURE
## LeNet

# FULLY CONVOLUTION NEURAL NETWORK IMAGE SEGMENTATION

# convnets perform classification

< 1 millisecond

"defect class1"

1000-dim vector

end-to-end learning

5

# lots of pixels, little time?

~1/10 second

???

end-to-end learning

6

# a classification network



convolution     fully connected

227 × 227    55 × 55    27 × 27    13 × 13

"defect class 1"

# end-to-end, pixels-to-pixels network

convolution



H × W  ·  H/4 × W/4  ·  H/8 × W/8  ·  H/16 × W/16  ·  H/32 × W/32  ·  H × W

conv, pool,
nonlinearity

upsampling

pixelwise
output + loss

8

# MRI image -> Left ventricle

## 2nd Data Science BOWL competition

# DATA SCIENCE BOWL 2017

## Predicting Lung Cancer



Lung segmentation → Nodule expert classification

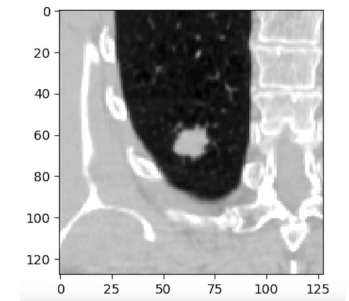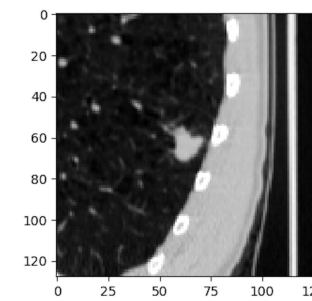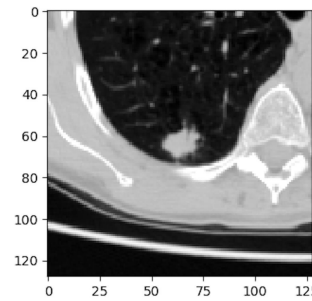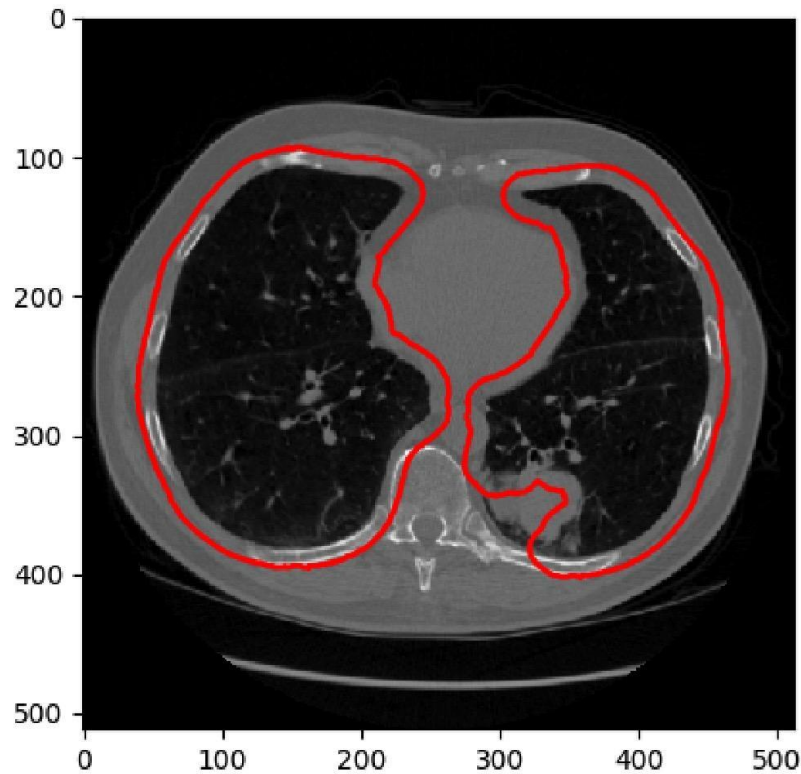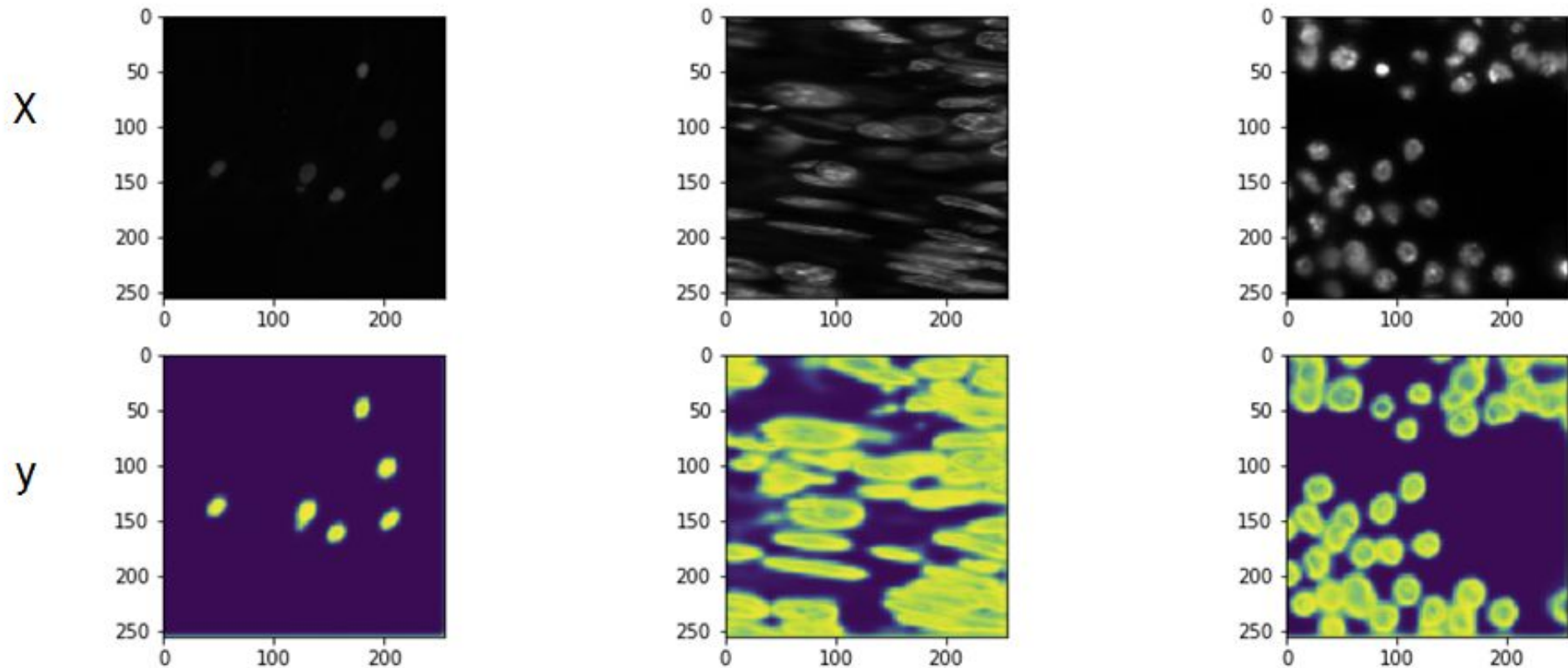# DATA SCIENCE BOWL 2018
## Predicting nuclei in divergent images

NVIDIA.

# INDUSTRIAL DEFECT INSPECTION

# DATASET

# INDUSTRIAL OPTICAL INSPECTION
## German Association for Pattern Recognition

# INDUSTRIAL OPTICAL INSPECTION
## German Association for Pattern Recognition



Pass

NG

Pass

NG

Pass

NG

Pass

NG

# DATA DETAILS

- Original images are 512 x 512 grayscale format

- Output is a tensor of size 512 x 512 x 1

  - Each pixel belongs to one of two classes

- Training set consist of 100 images

- Test set consist of 50 images

# MODEL SET UP

# Deconvolution layer

- Deconvolution (transpose convolution) layer

  - Up-sampling method to bring a smaller image data set back up to it's original size for final pixel classification

- Long et al (CVPR2015) has nice paper re: FCN for segmentation

  - Created FCNs from AlexNet and other canonical networks

- Zeiler et al (CVPR2010) describes deconvolution

# U-Net structure



https://arxiv.org/abs/1505.04597

3X3 Conv2d+ReLU

2X2 MaxPool

2X2 Conv2dTranspose

copy and concatenate

# IMBALANCE DATA

# Dice Metric

- Metric to compare the similarity of two samples:

$$\frac{2A_{nl}}{A_n + A_l}$$

IoU = $\dfrac{\text{Area of Overlap}}{\text{Area of Union}}$

  - Where:
    - $A_n$ is the area of the contour predicted by the network
    - $A_l$ is the area of the contour from the label
    - $A_{nl}$ is the intersection of the two
      - The area of the contour that is predicted correctly by the network
      - 1.0 means perfect score.

- More accurately compute how well we're predicting the contour against the label

- We can just count pixels to give us the respective areas

NVIDIA.

# APPLICATION: INDUSTRIAL INSPECTION
## NVIDIA

# FINAL DECISION
## Plus Human Logic



Size, Position, … etc

# PRODUCTION

# NVIDIA DEEP LEARNING SOFTWARE PLATFORM

## TRAINING

Training Data

Data Management

Training

Model Assessment

Trained Neural Network

Caffe2 · Chainer · Microsoft Cognitive Toolkit · mxnet · TensorFlow · PYTORCH · theano

## INFERENCE

Data center

GRE + TensorRT

Embedded

JETPACK SDK

Automotive

DriveWorks SDK

## NVIDIA DEEP LEARNING SDK and CUDA

cuDNN

NCCL

GPU0 · GPU1 · GPU3 · GPU2

cuBLAS

cuSPARSE

TensorRT

DeepStream SDK

NVDEC · Inference · NVENC

developer.nvidia.com/deep-learning-software

# CHALLENGES DURING PRODUCTION

| Requirement | Challenges |
|---|---|
| High Throughput | **Unable to processing high-volume, high-velocity data**<br>➢ Impact: Increased cost ($, time) per inference |
| Low Response Time | **Applications don't deliver real-time results**<br>➢ Impact: Negatively affects user experience (voice recognition, personalized recommendations, real-time object detection) |
| Power and Memory Efficiency | **Inefficient applications**<br>➢ Impact: Increased cost (running and cooling), makes deployment infeasible |
| Deployment-Grade Solution | **Research frameworks not designed for production**<br>➢ Impact: Framework overhead and dependencies increases time to solution and affects productivity |

# NVIDIA TENSORRT
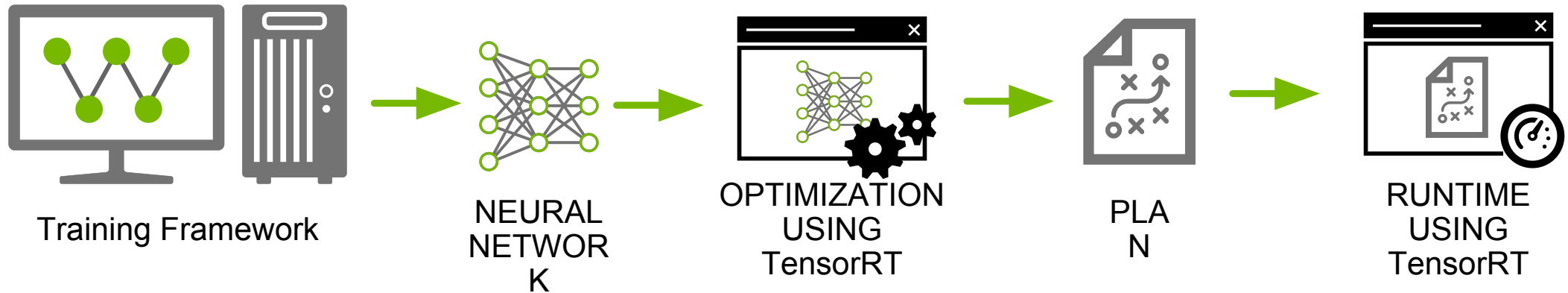## Programmable Inference Accelerator

# TENSORRT
## Workflow



Training Framework → NEURAL NETWORK → OPTIMIZATION USING TensorRT → PLAN → RUNTIME USING TensorRT
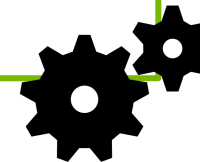
# TENSORRT
## Optimizations

**TRAINED NEURAL NETWORK**

- **Fuse network layers**
- **Eliminate concatenation layers**
- **Kernel specialization**
- **Auto-tuning for target platform**
- **Tuned for given batch size**

OPTIMIZED INFERENCE RUNTIME

# CHALLENGES ADDRESSED BY TENSORRT

| Requirement | TensorRT Delivers |
|---|---|
| **High Throughput** | **Maximizes inference performance on NVIDIA GPUs**<br>➤ INT8, FP16 Precision Calibration, Layer & Tensor Fusion, Kernel Auto-Tuning |
| **Low Response Time** | ➤ Up to 40x Faster than CPU-Only inference and 18x faster inference of TensorFlow models<br>➤ Under 7ms real-time latency |
| **Power and Memory Efficiency** | **Performs target specific optimizations**<br>➤ Platform specific kernels for Embedded (Jetson), Datacenter (Tesla GPUs) and Automotive (DrivePX)<br>➤ Dynamic Tensor Memory management improves memory re-use |
| **Deployment-Grade Solution** | **Designed for production environments**<br>➤ No framework overhead, minimal dependencies<br>➤ Multiple frameworks, Network Definition API<br>➤ C++, Python API, Customer Layer API |

 NVIDIA.

# THANKS!